

When is sparse dictionary learning well-posed?

Charles J. Garfinkle and Christopher J. Hillar
Redwood Center for Theoretical Neuroscience, Berkeley, CA, USA

Abstract—Dictionary learning methods for sparse coding have exposed underlying structure in many kinds of natural signals. However, universal theorems guaranteeing the statistical consistency of inference in this model are lacking. Here, we prove that for almost all diverse enough datasets generated from the model, latent dictionaries and sparse codes are uniquely identifiable up to an error commensurate with measurement noise. Applications are given to data analysis, neuroscience, and engineering.

I. INTRODUCTION

THE emergence of response properties of neurons in the mammalian visual cortex from the optimization of dictionaries for sparse coding of natural images marked an exciting development in computational neuroscience [1]–[4]. Many dictionary learning algorithms have since been developed and applied to a variety of problems in signal processing and machine learning (see [5] for a comprehensive review). A popular formulation of the idea is to encode each of N data points as a linear combination of at most k n -dimensional vectors from an inferred dictionary of size m , where $k < m \ll N$.

Certain applications to data analysis call for a unique such “sparse structure”. For instance, detecting forgeries by analysis of local painting style [6], [7] requires that all dictionaries consistent with training data do not differ appreciably in their ability to sparsely encode new samples. Recently, algorithms with proven convergence under certain conditions have been proposed (see [8, Sec. I–E] for a brief summary of the state-of-the-art). It has remained unknown, however, when the problem is well-posed in general (as per Hadamard [9]) – that is, independent of constraints specific to a particular algorithm.

The main finding of this work is that any dictionary satisfying the spark condition (2) from compressive sensing is identifiable from as few as $N = m + m(k-1)\binom{m}{k}$ noisy sparse linear combinations of its elements up to an error linear in the noise (Thm. 1). In fact, provided (n, m, k) satisfy the nearly-optimal compressive sensing inequality (6), in almost all cases the dictionary learning problem is well-posed given enough data (Cor. 2). Moreover, these guarantees extend easily to the case when only an upper bound on m is known (Thm. 3). We hope the explicit, algorithm-independent criteria we provide here may serve as a useful tool in future analyses of dictionary learning procedures.

We state the dictionary learning problem considered here more precisely as follows. Fix a dictionary represented as the columns A_j of a matrix $A \in \mathbb{R}^{n \times m}$ and suppose a dataset Z consists of measurements:

$$\mathbf{z}_i = A\mathbf{a}_i + \mathbf{n}_i, \quad i = 1, \dots, N, \quad (1)$$

for k -sparse $\mathbf{a}_i \in \mathbb{R}^m$ having at most k nonzero entries and noise $\mathbf{n}_i \in \mathbb{R}^n$, with bounded norm $\|\mathbf{n}_i\|_2 \leq \eta$ representing our combined worst-case uncertainty in measuring $A\mathbf{a}_i$.

Problem 1 (Dict. Learning). *Find $B \in \mathbb{R}^{n \times m}$ and k -sparse $\mathbf{b}_1, \dots, \mathbf{b}_N \in \mathbb{R}^m$ such that $\|\mathbf{z}_i - B\mathbf{b}_i\|_2 \leq \eta$ for $i = 1, \dots, N$.*

Note that any particular solution to this problem in fact represents a whole class of equivalent dictionaries BPD and codes $D^{-1}P^\top \mathbf{b}_i$, where $P \in \mathbb{R}^{m \times m}$ is any permutation matrix and $D \in \mathbb{R}^{m \times m}$ any invertible diagonal matrix. Since arbitrary scalings and ordering of dictionary elements represent the same underlying model, it is natural to ask whether solutions to Problem 1 are unique up to this equivalence.

Previous work [10]–[13] on the noiseless case $\eta = 0$ has shown that the solution (when it exists) is indeed unique in this sense provided the \mathbf{a}_i are sufficiently diverse and the matrix A satisfies the *spark condition*:

$$A\mathbf{x}_1 = A\mathbf{x}_2 \implies \mathbf{x}_1 = \mathbf{x}_2, \quad \text{for all } k\text{-sparse } \mathbf{x}_1, \mathbf{x}_2, \quad (2)$$

which is evidently a necessary condition for uniqueness given only that \mathbf{a}_i are k -sparse. Matrices of the form PD thus form the *ambiguity transformation group* inherent to the noiseless problem subject to these constraints [14].

We introduce the following terminology to handle $\eta > 0$.

Definition 1. Fix $Y = \{\mathbf{y}_1, \dots, \mathbf{y}_N\} \subset \mathbb{R}^n$. We say Y has a *k -sparse representation in \mathbb{R}^m* if there exists an $A \in \mathbb{R}^{n \times m}$ and k -sparse $\mathbf{a}_1, \dots, \mathbf{a}_N \in \mathbb{R}^m$ such that $\mathbf{y}_i = A\mathbf{a}_i$ for all i . This representation is *stable* if for every $\delta_1, \delta_2 \geq 0$, there exists $\varepsilon = \varepsilon(\delta_1, \delta_2) \geq 0$ (with $\varepsilon > 0$ when $\delta_1, \delta_2 > 0$) such that if a matrix $B \in \mathbb{R}^{n \times m}$ and k -sparse $\mathbf{b}_1, \dots, \mathbf{b}_N \in \mathbb{R}^m$ have $\|A\mathbf{a}_i - B\mathbf{b}_i\|_2 \leq \varepsilon$ for all i , then there is a permutation matrix $P \in \mathbb{R}^{m \times m}$ and invertible diagonal matrix $D \in \mathbb{R}^{m \times m}$ such that for all $i = 1, \dots, N$ and $j = 1, \dots, m$:

$$\|A_j - (BPD)_j\|_2 \leq \delta_1 \quad \text{and} \quad \|\mathbf{a}_i - D^{-1}P^\top \mathbf{b}_i\|_1 \leq \delta_2. \quad (3)$$

We ask here: *When does $Y \subset \mathbb{R}^n$ have a stable k -sparse representation in \mathbb{R}^m ?* To see how an affirmative answer to this question informs the interpretation of solutions to Problem 1, suppose that $Y = \{A\mathbf{a}_1, \dots, A\mathbf{a}_N\}$ has a stable k -sparse representation in \mathbb{R}^m , and fix δ_1, δ_2 to be the desired accuracy in recovery (3). Consider now any dataset Z generated as in (1) that has $\eta \leq \frac{1}{2}\varepsilon(\delta_1, \delta_2)$. Then, any dictionary B and k -sparse $\mathbf{b}_1, \dots, \mathbf{b}_N$ solving Problem 1 are necessarily close to the original matrix A and codes \mathbf{a}_i (i.e., satisfy (3)).

In the next section, we give precise statements of our main results, which include an explicit form for $\varepsilon(\delta_1, \delta_2)$. We then prove our main theorem (Thm. 1) in Sec. III after listing some additional definitions and lemmas required for the proof, including our main tool from combinatorial matrix analysis

(Lem. 1). Our proof is a refinement of the arguments in [13] to handle noise and to reduce the sufficient number of samples from $N = k \binom{m}{k}^2$ to $N = m + m(k-1) \binom{m}{k}$. All other proofs are relegated to the appendices. Finally, several applications are considered in our discussion, Sec. IV.

II. RESULTS

Before precisely stating our results, we explain how the spark condition (2) relates to the *lower bound* [15] of A , written $L(A)$, which is the largest number α such that $\|Ax\|_2 \geq \alpha \|x\|_2$ for all $x \in \mathbb{R}^m$. By compactness, every injective linear map has a nonzero lower bound; hence, if A satisfies (2), then every submatrix formed from $2k$ of its columns or less has a nonzero lower bound. We therefore define the following domain-restricted lower bound of A :

$$L_k(A) := \max\{\alpha : \|Ax\|_2 \geq \alpha \|x\|_2 \text{ for all } k\text{-sparse } x \in \mathbb{R}^m\}.$$

Clearly, $L_k(A) \geq L_{k'}(A)$ whenever $k < k'$, and for any A satisfying (2), we have $L_{k'}(A) > 0$ for all $k' \leq 2k$.

A *cyclic order* on $[m] := \{1, \dots, m\}$ is an arrangement of $[m]$ in a circular necklace, and an *interval* in the order is any subset of contiguous elements. A vector $\mathbf{a} \in \mathbb{R}^m$ is said to be *supported* on $S \subseteq [m]$ when $\mathbf{a} \in \text{Span}_{\mathbb{R}}\{\mathbf{e}_j\}_{j \in S}$, where \mathbf{e}_j are the standard basis vectors. Also, recall that M_j denotes the j th column of a matrix M . The following result gives an answer to our question from the introduction.

Theorem 1. Fix n, m , and $k < m$. If $A \in \mathbb{R}^{n \times m}$ is injective on k -sparse vectors and k -sparse $\mathbf{a}_1, \dots, \mathbf{a}_N \in \mathbb{R}^m$ are such that for every interval of length k in some cyclic order on $[m]$ there are $(k-1) \binom{m}{k} + 1$ vectors \mathbf{a}_i in general linear position (i.e., any k of them are linearly independent) supported there, then $Y = \{A\mathbf{a}_1, \dots, A\mathbf{a}_N\}$ has a stable k -sparse representation in \mathbb{R}^m .

Specifically, there exists a constant $C > 0$ for which the following holds for all $\varepsilon < \frac{L_2(A)}{\sqrt{2}} C^{-1}$.¹ If any matrix $B \in \mathbb{R}^{n \times m}$ and k -sparse $\mathbf{b}_1, \dots, \mathbf{b}_N \in \mathbb{R}^m$ are such that $\|A\mathbf{a}_i - B\mathbf{b}_i\|_2 \leq \varepsilon$ for all $i \in [N]$, then for all $j \in [m]$:

$$\|A_j - (BPD)_j\|_2 \leq C\varepsilon, \quad (4)$$

for some permutation matrix $P \in \mathbb{R}^{m \times m}$ and invertible diagonal matrix $D \in \mathbb{R}^{m \times m}$. Moreover, if $\varepsilon < \varepsilon_0 := \frac{L_{2k}(A)}{\sqrt{2k}} C^{-1}$, then $L_{2k}(BPD) \geq L_{2k}(A) (1 - \varepsilon/\varepsilon_0)$ and for all $i \in [N]$:

$$\|\mathbf{a}_i - D^{-1}P^\top \mathbf{b}_i\|_1 \leq \frac{\varepsilon}{\varepsilon_0 - \varepsilon} (C^{-1} + \|\mathbf{a}_i\|_1). \quad (5)$$

The constant C is explicitly defined in (7) below.

An important consequence of this result is that (3) is guaranteed provided ε does not exceed:

$$\varepsilon(\delta_1, \delta_2) := \min \left\{ \frac{\delta_1}{C}, \frac{\delta_2 L_{2k}(A)/\sqrt{2k}}{1 + C (\max_{i \in [N]} \|\mathbf{a}_i\|_1 + \delta_2)} \right\}.$$

Corollary 1. Given n, m , and $k < m$, there are $N = m(k-1) \binom{m}{k} + m$ vectors $\mathbf{a}_1, \dots, \mathbf{a}_N \in \mathbb{R}^m$ such that every

¹The condition $\varepsilon < \frac{L_2(A)}{\sqrt{2}} C^{-1}$ is necessary; otherwise, with $A = I$ and $\mathbf{a}_i = \mathbf{e}_i$, there is a B and 1-sparse \mathbf{b}_i with $\|A\mathbf{a}_i - B\mathbf{b}_i\|_2 \leq \varepsilon$ violating (4).

matrix $A \in \mathbb{R}^{n \times m}$ satisfying (2) generates a set $Y = \{A\mathbf{a}_1, \dots, A\mathbf{a}_N\}$ with a stable k -sparse representation in \mathbb{R}^m .

It is straightforward to provide a probabilistic extension of Thm. 1 using the following fact in random matrix theory. The matrix $A \in \mathbb{R}^{n \times m}$ satisfies (2) with probability one (or with “high probability” for discrete variables) provided:

$$n \geq \gamma k \log \left(\frac{m}{k} \right), \quad (6)$$

where γ is a positive constant dependent on the particular distribution from which the entries of A are sampled i.i.d. (many ensembles suffice, e.g. [16, Sec. 4]).

In fact, the spark condition can be made explicit. Let A be the $n \times m$ matrix of nm indeterminates A_{ij} . When real numbers are substituted for all the A_{ij} , the resulting matrix satisfies (2) if and only if the following polynomial is nonzero:

$$f(A) = \prod_{S \in \binom{[m]}{k}} \sum_{S' \in \binom{[n]}{k}} (\det A_{S', S})^2,$$

where for any $S' \in \binom{[n]}{k}$ and $S \in \binom{[m]}{k}$, the symbol $A_{S', S}$ denotes the submatrix of entries A_{ij} with $(i, j) \in S' \times S$.

Since f is a real analytic function, it is enough to show that at least *one* substitution of real numbers satisfies $f(A) \neq 0$ to conclude that its zeroes form a set with measure zero. Hence, every $n \times m$ matrix A outside a set of measure zero satisfies (2) provided (6) holds for a value of γ for *some* distribution. We remark that the precise relationship between m , n , and k guaranteeing that f is not identically zero is still an open problem in real algebraic geometry. In any case, we set γ_0 to be the smallest known such γ .

It so happens that a similar statement applies to sets of vectors with a stable sparse representation. As in [13, Sec. IV], consider the “symbolic” dataset $Y = \{A\mathbf{a}_1, \dots, A\mathbf{a}_N\}$ generated by indeterminate A and k -sparse $\mathbf{a}_1, \dots, \mathbf{a}_N$.

Theorem 2. Fix n, m , and $k < m$. There is a polynomial in the entries of A and the \mathbf{a}_i with the following property: if the polynomial evaluates to a nonzero number and for every interval of length k in some cyclic order on $[m]$ at least $(k-1) \binom{m}{k} + 1$ of the resulting vectors \mathbf{a}_i are supported on that interval, then Y has a stable k -sparse representation in \mathbb{R}^m (Def. 1). In particular, either no substitutions impart to Y this property or all but a Borel set of measure zero do.

Corollary 2. Fix n, m , and k satisfying (6) for $\gamma = \gamma_0$, and let the entries of the matrix $A \in \mathbb{R}^{n \times m}$ and k -sparse vectors $\mathbf{a}_1, \dots, \mathbf{a}_N \in \mathbb{R}^m$ be drawn independently from probability measures absolutely continuous with respect to the standard Borel measure μ . If at least $(k-1) \binom{m}{k} + 1$ of the vectors \mathbf{a}_i are supported on each interval of length k in some cyclic order on $[m]$, then Y has a stable k -sparse representation in \mathbb{R}^m with probability one.

Next, we address when only an upper bound m' on the latent dimension m is known (assuming that B satisfies (2)).

Theorem 3. Let Y be defined as in Thm. 1. There exists $C > 0$ for which the following holds for all $\varepsilon < \frac{L_2(A)}{\sqrt{2}} C^{-1}$ and any $m' > m$. If a matrix $B \in \mathbb{R}^{n \times m'}$ satisfies (2) and k -sparse

$\mathbf{b}_1, \dots, \mathbf{b}_N \in \mathbb{R}^{m'}$ are such that $\|\mathbf{A}\mathbf{a}_i - \mathbf{B}\mathbf{b}_i\|_2 \leq \varepsilon$ for all $i \in [N]$, then (4) and (5) hold for some $n \times m$ submatrix of B and corresponding subvectors of the \mathbf{b}_i , respectively.

In other words, the columns of B contain (up to noise, after appropriate scaling) the columns of the original dictionary A . Similarly, the \mathbf{b}_i contain the original codes \mathbf{a}_i . The constant C here is expression (22) from the proof of Thm. 3.

III. PROOF OF THEOREM 1

We first briefly outline our main tools, which include general notions of angle (Def. 2) and distance (Def. 4) between subspaces as well as a (stable) uniqueness result in matrix analysis (Lem. 1). Given a set \mathcal{T} , let $\binom{\mathcal{T}}{k}$ denote the set of all subsets of \mathcal{T} of size k . Denote by \mathbf{e}_i for $i \in [m]$ the canonical basis vectors in \mathbb{R}^m and let $\text{Span}\{\mathbf{v}_1, \dots, \mathbf{v}_\ell\}$ be the \mathbb{R} -linear span of vectors $\mathbf{v}_1, \dots, \mathbf{v}_\ell$. Given $S \subseteq [m]$ and $M \in \mathbb{R}^{n \times m}$, let M_S be the submatrix with columns M_j for $j \in S$, which will also denote its column space when appropriate. Let \mathfrak{S}_k denote the set of permutations (bijections) on k elements.

Definition 2. The *angle* between subspaces $U, V \subseteq \mathbb{R}^n$ is the angle $\theta(U, V)$ in $(0, \frac{\pi}{2}]$ whose cosine is [17, Def. 9.4]:

$$\cos \theta(U, V) := \max \left\{ |\langle \mathbf{u}, \mathbf{v} \rangle| : \begin{array}{l} \mathbf{u} \in U \cap (U \cap V)^\perp \cap \mathcal{B} \\ \mathbf{v} \in V \cap (U \cap V)^\perp \cap \mathcal{B} \end{array} \right\},$$

where $\mathcal{B} = \{\mathbf{x} \in \mathbb{R}^n : \|\mathbf{x}\|_2 \leq 1\}$.

The next quantity is based on one used in [17] to analyze the convergence of the alternating projections algorithm for projecting a point onto the intersection of a set of subspaces.

Definition 3. Fix $A \in \mathbb{R}^{n \times m}$ and $k < m$ and let $\mathcal{T} \subseteq \binom{[m]}{k}$ be of size at least k . For $k > 1$, define:

$$\phi_{\mathcal{T}}(A) := \min_{T \in \binom{\mathcal{T}}{k}} 1 - \xi(\{A_S : S \in T\}),$$

where for any set $\mathcal{V} = \{V_1, \dots, V_k\}$ of subspaces of \mathbb{R}^m :

$$\xi(\mathcal{V}) := \min_{\sigma \in \mathfrak{S}_k} \left(1 - \prod_{i=1}^{k-1} \sin^2 \theta(V_{\sigma(i)}, \cap_{j>i} V_{\sigma(j)}) \right)^{1/2},$$

and we set $\phi_{\mathcal{T}}(A) := 1$ for $k = 1$.

Note that $\xi < 1$ and $\phi_{\mathcal{T}}(A) > 0$ always, since it is clear from Def. 2 that $\cos \theta(U, V) < 1$ for all subspaces $U, V \in \mathbb{R}^m$.

We now state the constant $C > 0$ referred to in Thm. 1:

$$C = \left(\frac{\sqrt{k^3}}{\phi_{\mathcal{T}}(A)} \right) \frac{\max_{j \in [m]} |A_j|_2}{\min_{S \in \mathcal{T}} L_k(AX_{I(S)})}, \quad (7)$$

where \mathcal{T} is the set intervals in the cyclic order, X is the $m \times N$ matrix with columns \mathbf{a}_i , and $I(S) := \{i : \text{supp}(\mathbf{a}_i) = S\}$.²

Definition 4. The *gap* (or *aperture*) $\Theta(U, V)$ between subspaces $U, V \subseteq \mathbb{R}^m$ is given by [18, Sec. 2]:

$$\Theta(U, V) := \max \left\{ \sup_{\mathbf{u} \in U, \|\mathbf{u}\|_2=1} d(\mathbf{u}, V), \sup_{\mathbf{v} \in V, \|\mathbf{v}\|_2=1} d(\mathbf{v}, U) \right\},$$

where $d(\mathbf{u}, V) := \inf\{\|\mathbf{u} - \mathbf{v}\|_2 : \mathbf{v} \in V\}$.

²That $C > 0$ is well-defined follows from $L_k(A) > 0$ and the general linear position of the \mathbf{a}_i , since then $|A_j|_2 > 0$ for all j and $\min_{S \in \mathcal{T}} L_k(AX_{I(S)}) > 0$.

We now state our uniqueness result in matrix analysis, generalizing [13, Lem. 1] to the noisy case.

Lemma 1 (Main Lemma). Fix $n, m, k < m$, and let \mathcal{T} be the set of intervals of length k in some cyclic order on $[m]$. Let $A, B \in \mathbb{R}^{n \times m}$ and suppose that A satisfies the spark condition (2) and has maximum column ℓ_2 -norm ρ . If there exists a map $\pi : \mathcal{T} \rightarrow \binom{[m]}{k}$ and some $\delta < \frac{L_2(A)}{\sqrt{2}}$ such that:

$$\Theta(A_S, B_{\pi(S)}) \leq \frac{\phi_{\mathcal{T}}(A)}{\rho k} \delta, \quad \text{for all } S \in \mathcal{T}, \quad (8)$$

then there exist a permutation matrix $P \in \mathbb{R}^{m \times m}$ and an invertible diagonal matrix $D \in \mathbb{R}^{m \times m}$ with:

$$|A_j - (BPD)_j|_2 \leq \delta, \quad \text{for } j \in [m]. \quad (9)$$

We will also use a few facts about d in Def. 4. The first, proven in [19, Lem. 3.2], is that if $\dim(W) = \dim(V)$ then:

$$\sup_{\mathbf{v} \in V, \|\mathbf{v}\|_2=1} d(\mathbf{v}, W) = \sup_{\mathbf{w} \in W, \|\mathbf{w}\|_2=1} d(\mathbf{w}, V). \quad (10)$$

The second is:

Lemma 2. If U, V are subspaces of \mathbb{R}^m , then:

$$d(\mathbf{u}, V) < \|\mathbf{u}\|_2 \text{ and } \mathbf{u} \in U \setminus \{\mathbf{0}\} \implies \dim(U) \leq \dim(V).$$

Finally, we often use $\|\mathbf{x}\|_1 \leq \sqrt{k} \|\mathbf{x}\|_2$ for k -sparse $\mathbf{x} \in \mathbb{R}^m$.

Let us now prove Thm. 1 for the simple case when $k = 1$. Fix $A \in \mathbb{R}^{n \times m}$ satisfying (2), and let $\mathbf{a}_i = c_i \mathbf{e}_i$ for $c_i \in \mathbb{R} \setminus \{0\}$, $i \in [m]$. By (7), we have:

$$C = \sqrt{k^3} \left(\frac{\max_{j \in [m]} |A_j|_2}{\min_{j \in [m]} |c_j A_j|_2} \right) \geq \left(\min_{\ell \in [m]} |c_\ell| \right)^{-1}. \quad (11)$$

Suppose that for some $B \in \mathbb{R}^{n \times m}$ and 1-sparse $\mathbf{b}_i \in \mathbb{R}^m$ we have $\|\mathbf{A}\mathbf{a}_i - \mathbf{B}\mathbf{b}_i\|_2 \leq \varepsilon < \frac{L_2(A)}{\sqrt{2}} C^{-1}$ for $i \in [m]$. Then there are $c'_1, \dots, c'_m \in \mathbb{R}$ and $\pi : [m] \rightarrow [m]$ with:

$$|c_j A_j - c'_j B_{\pi(j)}|_2 \leq \varepsilon, \quad \text{for } j \in [m]. \quad (12)$$

Note that $c'_j \neq 0$ for all j since otherwise (by definition of $L_2(A)$), we would have $|c_j A_j|_2 < \min_{\ell \in [m]} |c_\ell A_\ell|_2$.

We now show that π is injective (and thus is a permutation). Suppose that $\pi(i) = \pi(j) = \ell$ for some $i \neq j$ and ℓ . Then, $|c_j A_j - c'_j B_\ell|_2 \leq \varepsilon$ and $|c_i A_i - c'_i B_\ell|_2 \leq \varepsilon$. Scaling and summing these inequalities by $|c'_i|$ and $|c'_j|$, respectively, and applying the triangle inequality, we obtain:

$$\begin{aligned} (|c'_i| + |c'_j|)\varepsilon &\geq |A(c'_i c_j \mathbf{e}_j - c'_j c_i \mathbf{e}_i)|_2 \\ &\geq \frac{L_2(A)}{\sqrt{2}} (|c'_i| + |c'_j|) \min_{\ell \in [m]} |c_\ell|. \end{aligned} \quad (13)$$

Since (13) contradicts (11) and our upper bound on ε , the map π is injective. Letting $P = (\mathbf{e}_{\pi(1)} \cdots \mathbf{e}_{\pi(m)})$ and $D = \text{diag}(c'_1, \dots, c'_m)$, we see that (12) becomes (4) for $j \in [m]$:

$$|A_j - (BPD)_j|_2 = |A_j - \frac{c'_j}{c_j} B_{\pi(j)}|_2 \leq \frac{\varepsilon}{|c_j|} \leq C\varepsilon.$$

It turns out that (4) already implies in general the recovery result (5) for $k \geq 1$ when $\varepsilon < \varepsilon_0 := \frac{L_2(A)}{2k} C^{-1}$. To see why, note that for all $2k$ -sparse $\mathbf{x} \in \mathbb{R}^m$, the triangle inequality gives $\|(A - BPD)\mathbf{x}\|_2 \leq C\varepsilon \|\mathbf{x}\|_1 \leq C\varepsilon \sqrt{2k} \|\mathbf{x}\|_2$. Thus,

$|BPD\mathbf{x}|_2 \geq ||A\mathbf{x}|_2 - |(A - BPD)\mathbf{x}|_2| \geq (L_{2k}(A) - \sqrt{2k}C\varepsilon)|\mathbf{x}|_2$, where we drop the absolute value since $\varepsilon < \varepsilon_0$. Hence, $L_{2k}(BPD) \geq L_{2k}(A)(1 - \varepsilon/\varepsilon_0) > 0$ and (5) then follows from:

$$\begin{aligned} |D^{-1}P^\top \mathbf{b}_i - \mathbf{a}_i|_1 &\leq \frac{\sqrt{2k}}{L_{2k}(BPD)} |BPD(\mathbf{a}_i - D^{-1}P^\top \mathbf{b}_i)|_2 \\ &\leq \frac{\varepsilon\sqrt{2k}}{L_{2k}(BPD)} (1 + C|\mathbf{a}_i|_1). \end{aligned}$$

It remains to show that (4) with C given in (7) follows from $\varepsilon < \frac{L_2(A)}{\sqrt{2}}C^{-1}$ for $k > 1$. Our main tool is Lem. 1.

Proof of Thm. 1: From above, we may assume that $k > 1$. Let \mathcal{T} be the set of intervals of length k in the given cyclic order on $[m]$, and fix $A \in \mathbb{R}^{n \times m}$ satisfying (2) and $N = m(k-1)\binom{m}{k} + m$ vectors $\mathbf{a}_i \in \mathbb{R}^m$ as in the statement of the theorem.

Suppose that for some $B \in \mathbb{R}^{n \times m}$ there exist k -sparse $\mathbf{b}_i \in \mathbb{R}^m$ such that $|A\mathbf{a}_i - B\mathbf{b}_i|_2 \leq \varepsilon$ for all $i \in [N]$. Since there are $(k-1)\binom{m}{k} + 1$ vectors \mathbf{a}_i with a given support $S \in \mathcal{T}$, the pigeon-hole principle implies that there exists some $S' \in \binom{[m]}{k}$ and some set of k indices $J(S)$ such that all \mathbf{a}_i and \mathbf{b}_i with $i \in J(S)$ have supports S and S' , respectively.

Let X and X' be the $m \times N$ matrices with columns \mathbf{a}_i and \mathbf{b}_i , respectively. It follows from the general linear position of the \mathbf{a}_i and the linear independence of every k columns of A that $L(AX_{J(S)}) > 0$; that is, the columns of the $n \times k$ matrix $AX_{J(S)}$ are linearly independent and thus form a basis for $\text{Span}\{A_S\}$. Fixing $\mathbf{y} \in \text{Span}\{A_S\}$, there then exists a unique $\mathbf{c} = (c_1, \dots, c_k) \in \mathbb{R}^k$ such that $\mathbf{y} = AX_{J(S)}\mathbf{c}$. Letting $\mathbf{y}' = BX'_{J(S)}\mathbf{c}$, which is in $\text{Span}\{B_{S'}\}$, we have:

$$\begin{aligned} |\mathbf{y} - \mathbf{y}'|_2 &= \left| \sum_{i=1}^k c_i (AX_{J(S)} - BX'_{J(S)})_i \right|_2 \leq \varepsilon \sum_{i=1}^k |c_i| \\ &\leq \varepsilon\sqrt{k}|\mathbf{c}|_2 \leq \frac{\varepsilon\sqrt{k}}{L(AX_{J(S)})} |AX_{J(S)}\mathbf{c}|_2 = \frac{\varepsilon\sqrt{k}}{L(AX_{J(S)})} |\mathbf{y}|_2. \end{aligned}$$

In particular, the inequality $d(\mathbf{y}, B_{S'}) \leq \varepsilon\sqrt{k}/L(AX_{J(S)})$ holds for all $\mathbf{y} \in \text{Span}\{A_S\}$ having unit ℓ_2 -norm.

We now show that (4) follows if $\varepsilon < \frac{L_2(A)}{\sqrt{2}}C^{-1}$, for C as defined in (7). Letting $\rho = \max_{j \in [m]} |A_j|_2$, we have:

$$\sup_{\substack{\mathbf{y} \in \text{Span}\{A_S\} \\ |\mathbf{y}|_2=1}} d(\mathbf{y}, B_{S'}) \leq \frac{\varepsilon\sqrt{k}}{L(AX_{J(S)})} < \frac{\phi_{\mathcal{T}}(A)L_2(A)}{\rho k\sqrt{2}}. \quad (14)$$

Since $L_2(A) \leq \rho\sqrt{2}$ and $\phi_{\mathcal{T}}(A) \leq 1$, the RHS of (14) is strictly less than one. It follows by Lem. 2 that $\dim(\text{Span}\{B_{S'}\}) \geq \dim(\text{Span}\{A_S\}) = k$. Since $|S'| = k$, we also have $\dim(\text{Span}\{B_{S'}\}) \leq k$; hence, $\dim(\text{Span}\{B_{S'}\}) = \dim(\text{Span}\{A_S\})$. Recalling (10), we see that the association $S \mapsto S'$ therefore defines a map $\pi : \mathcal{T} \rightarrow \binom{[m]}{k}$ satisfying:

$$\Theta(A_S, B_{\pi(S)}) \leq \frac{\varepsilon\sqrt{k}}{L(AX_{J(S)})}, \quad \text{for } S \in \mathcal{T}. \quad (15)$$

Thus, from (14) and (15) it follows that, for any particular $S \in \mathcal{T}$, the inequality in (8) is satisfied for:

$$\delta = \frac{\rho k}{\phi_{\mathcal{T}}(A)} \left(\frac{\varepsilon\sqrt{k}}{L(AX_{J(S)})} \right) < \frac{L_2(A)}{\sqrt{2}}.$$

Hence, (8) holds (for all $S \in \mathcal{T}$) when $\delta = C\varepsilon$. The result (4) follows by application of Lem. 1 and, as demonstrated previously, (4) implies (5). ■

IV. DISCUSSION

In this note, we generalize recent results [13] on the uniqueness of solutions to Problem 1 to the case of noisy measurements while also significantly reducing the number of required samples. We remark that our result in the deterministic case (Thm. 1) accounts for *worst-case* noise, whereas the “effective” noise might be much smaller when it is sampled from a given distribution; in such cases, the constants C in Thms. 1, 3 may be much smaller with high probability. We note also that these results extend trivially to when point-wise injective nonlinearities are applied to the data. We close by outlining four diverse application areas.

Inverse Problems. Our results provide theoretical grounding for the use of sparse dictionary learning in blind source separation, wherein the goal is to infer the generating dictionary and sparse codes from noisy measurements generated as in (1) (e.g., recovering a rat’s position on a linear track from local field potentials in Hippocampus [20]). It would be of practical utility therefore to determine the best possible dependence of ε on δ_1, δ_2 (see Thm. 1) as well as the minimal requirements on the number and diversity of generating codes.

Theoretical Neuroscience. Sparse dictionary learning and related methods have recovered characteristic components of natural images [1]–[4] and sounds [21]–[23], reproducing response properties of cortical neurons. Our theorems suggest that this correspondence could be due to the uniqueness of sparse representations. Furthermore, our guarantees justify the hypothesis of [24] and [25] that sparse codes passed through a communication bottleneck in the brain can be recovered from random projections via (unsupervised) biologically plausible sparse dictionary learning (e.g., [26], [27]).

Smoothed Analysis. The main concept in smoothed analysis [28] is that certain algorithms having exponential worst-case behavior are, nonetheless, efficient if certain (typically, measure zero in the continuous case and with “low probability” in the discrete case) pathological input sets are avoided. Our results imply that if there is an efficient “smoothed” algorithm for solving Problem 1 given enough samples, then for generic inputs this algorithm determines the unique original solution. We note that avoiding “bad” (NP-hard) sets of inputs is a necessary technicality for dictionary learning [29], [30].

Engineering. Several groups utilize compressive sensing for signal processing tasks: MRI analysis [31], image compression [32], and, more recently, the design of an ultrafast camera [33]. Given such effective uses of compressive sensing, it is only a matter of time before these systems incorporate sparse dictionary learning to encode and process data. Guarantees such as those offered by our theorems allow any such device to be equivalent to any other (having different initial parameters and data samples) as long as enough data originate from a statistically identical system.

We thank Fritz Sommer for turning our attention to the dictionary learning problem, Darren Rhea for sharing early explorations, and Ian Morris for posting identity (10) online.

REFERENCES

- [1] B. Olshausen and D. Field, "Emergence of simple-cell receptive field properties by learning a sparse code for natural images," *Nature*, vol. 381, no. 6583, pp. 607–609, 1996.
- [2] J. Hurri, A. Hyvärinen, J. Karhunen, and E. Oja, "Image feature extraction using independent component analysis," in *Proc. NORSIG '96 (Nordic Signal Processing Symposium)*, 1996, pp. 475–478.
- [3] A. Bell and T. Sejnowski, "The "independent components" of natural scenes are edge filters," *Vision Research*, vol. 37, no. 23, pp. 3327–3338, 1997.
- [4] J. van Hateren and A. van der Schaaf, "Independent component filters of natural images compared with simple cells in primary visual cortex," *Proceedings of the Royal Society of London. Series B: Biological Sciences*, vol. 265, no. 1394, pp. 359–366, 1998.
- [5] Z. Zhang, Y. Xu, J. Yang, X. Li, and D. Zhang, "A survey of sparse representation: algorithms and applications," *Access, IEEE*, vol. 3, pp. 490–530, 2015.
- [6] J. Hughes, D. Graham, and D. Rockmore, "Quantification of artistic style through sparse coding analysis in the drawings of Pieter Bruegel the Elder," *Proc. of the National Academy of Sciences*, vol. 107, no. 4, pp. 1279–1283, 2010.
- [7] B. Olshausen and M. DeWeese, "Applied mathematics: The statistics of style," *Nature*, vol. 463, no. 7284, pp. 1027–1028, 2010.
- [8] J. Sun, Q. Qu, and J. Wright, "Complete dictionary recovery over the sphere I: Overview and the geometric picture," *Information Theory, IEEE Transactions on*, 2016.
- [9] J. Hadamard, "Sur les problèmes aux dérivées partielles et leur signification physique," *Princeton University Bulletin*, vol. 13, no. 49-52, p. 28, 1902.
- [10] Y. Li, A. Cichocki, and S.-I. Amari, "Analysis of sparse representation and blind source separation," *Neural Computation*, vol. 16, no. 6, pp. 1193–1234, 2004.
- [11] P. Georgiev, F. Theis, and A. Cichocki, "Sparse component analysis and blind source separation of underdetermined mixtures," *IEEE Transactions on Neural Networks*, vol. 16, pp. 992–996, 2005.
- [12] M. Aharon, M. Elad, and A. Bruckstein, "On the uniqueness of overcomplete dictionaries, and a practical way to retrieve them," *Linear Algebra and its Applications*, vol. 416, no. 1, pp. 48–67, 2006.
- [13] C. Hillar and F. Sommer, "When can dictionary learning uniquely recover sparse data from subsamples?" *Information Theory, IEEE Transactions on*, vol. 61, no. 11, pp. 6290–6297, 2015.
- [14] Y. Li, K. Lee, and Y. Bresler, "A unified framework for identifiability analysis in bilinear inverse problems with applications to subspace and sparsity models," *arXiv preprint arXiv:1501.06120*, 2015.
- [15] J. Grcar, "A matrix lower bound," *Linear Algebra and its Applications*, vol. 433, no. 1, pp. 203–220, 2010.
- [16] R. Baraniuk, M. Davenport, R. DeVore, and M. Wakin, "A simple proof of the restricted isometry property for random matrices," *Constructive Approximation*, vol. 28, no. 3, pp. 253–263, 2008.
- [17] F. Deutsch, *Best approximation in inner product spaces*. Springer Science & Business Media, 2012.
- [18] T. Kato, *Perturbation theory for linear operators*. Springer Science & Business Media, 2013, vol. 132.
- [19] I. Morris, "A rapidly-converging lower bound for the joint spectral radius via multiplicative ergodic theory," *Advances in Mathematics*, vol. 225, no. 6, pp. 3425–3445, 2010.
- [20] G. Agarwal, I. Stevenson, A. Berényi, K. Mizuseki, G. Buzsáki, and F. Sommer, "Spatially distributed local fields in the hippocampus encode rat position," *Science*, vol. 344, no. 6184, pp. 626–630, 2014.
- [21] A. Bell and T. Sejnowski, "Learning the higher-order structure of a natural sound," *Network: Computation in Neural Systems*, vol. 7, no. 2, pp. 261–266, 1996.
- [22] E. Smith and M. Lewicki, "Efficient auditory coding," *Nature*, vol. 439, no. 7079, pp. 978–982, 2006.
- [23] N. Carlson, V. Ming, and M. DeWeese, "Sparse codes for speech predict spectrotemporal receptive fields in the inferior colliculus," *PLoS Comput Biol*, vol. 8, no. 7, p. e1002594, 2012.
- [24] W. Coulter, C. Hillar, G. Isley, and F. Sommer, "Adaptive compressed sensing – a new class of self-organizing coding models for neuroscience," in *Acoustics Speech and Signal Processing (ICASSP), 2010 IEEE International Conference on*. IEEE, 2010, pp. 5494–5497.
- [25] G. Isely, C. Hillar, and F. Sommer, "Deciphering subsampled data: adaptive compressive sampling as a principle of brain communication," in *Advances in Neural Information Processing Systems*, 2010, pp. 910–918.
- [26] C. Rozell, D. Johnson, R. Baraniuk, and B. Olshausen, "Neurally plausible sparse coding via thresholding and local competition," *Neural Computation*, 2007.
- [27] T. Hu, C. Pehlevan, and D. B. Chklovskii, "A Hebbian/anti-Hebbian network for online sparse dictionary learning derived from symmetric matrix factorization," in *2014 48th Asilomar Conference on Signals, Systems and Computers*. IEEE, 2014, pp. 613–619.
- [28] D. Spielman and S.-H. Teng, "Smoothed analysis of algorithms: Why the simplex algorithm usually takes polynomial time," *Journal of the ACM (JACM)*, vol. 51, no. 3, pp. 385–463, 2004.
- [29] M. Razaviyayn, H.-W. Tseng, and Z.-Q. Luo, "Computational intractability of dictionary learning for sparse representation," *arXiv preprint arXiv:1511.01776*, 2015.
- [30] A. Tillmann, "On the computational intractability of exact and approximate dictionary learning," *Signal Processing Letters, IEEE*, vol. 22, no. 1, pp. 45–49, 2015.
- [31] M. Lustig, D. Donoho, J. Santos, and J. Pauly, "Compressed sensing MRI," *IEEE Signal Processing Magazine*, vol. 25, no. 2, pp. 72–82, 2008.
- [32] M. Duarte, M. Davenport, D. Takbar, J. Laska, T. Sun, K. Kelly, and R. Baraniuk, "Single-pixel imaging via compressive sampling," *IEEE Signal Processing Magazine*, vol. 25, no. 2, pp. 83–91, March 2008.
- [33] L. Gao, J. Liang, C. Li, and L. Wang, "Single-shot compressed ultrafast photography at one hundred billion frames per second," *Nature*, vol. 516, no. 7529, pp. 74–77, 2014.
- [34] G. Folland, *Real analysis: modern techniques and their applications*. John Wiley & Sons, 2013.

Charles J. Garfinkle : B.S. Physics, B.S. Chemistry, McGill University, 2010. Currently, Ph.D. candidate, Neuroscience, University of California, Berkeley.

Christopher J. Hillar: B.S. Mathematics, B.S. Computer Science, Yale University, 2000; Ph.D., Mathematics, University of California, Berkeley, CA 2005. Currently, Redwood Center for Theoretical Neuroscience, UCB.

APPENDIX A
COMBINATORIAL MATRIX ANALYSIS

Here, we prove Lem. 1, which is the main ingredient in our proof of Thm. 1. We then outline how additionally assuming the spark condition (2) for B simplifies the proof and also allows for its extension to the case where only an upper bound on the number of columns m of A is known (Thm. 3).

We first prove some auxiliary lemmas, starting with a proof of Lem. 2, which was stated in Section III.

Proof of Lemma 2: We prove the contrapositive. If $\dim(U) > \dim(V)$, then a dimension argument ($\dim U + \dim V^\perp > m$) gives a nonzero $\mathbf{u} \in U \cap V^\perp$. In particular, we have $\|\mathbf{u} - \mathbf{v}\|_2^2 = \|\mathbf{u}\|_2^2 + \|\mathbf{v}\|_2^2 \geq \|\mathbf{u}\|_2^2$ for $\mathbf{v} \in V$, and thus $d(\mathbf{u}, V) \geq \|\mathbf{u}\|_2$. ■

Given sets \mathcal{T} , let $\cap \mathcal{T}$ denote their intersection.

Lemma 3. *Let $M \in \mathbb{R}^{n \times m}$. If every $2k$ columns of M are linearly independent, then for any $\mathcal{T} \subseteq \bigcup_{\ell \leq k} \binom{[m]}{\ell}$, we have:*

$$\text{Span}\{M_{\cap \mathcal{T}}\} = \bigcap_{S \in \mathcal{T}} \text{Span}\{M_S\}.$$

Proof: By induction, it is enough to prove the lemma when $|\mathcal{T}| = 2$. The proof now follows directly from the assumption. ■

Lemma 4. *Fix $k \geq 2$. Let $\mathcal{V} = \{V_1, \dots, V_k\}$ be subspaces of \mathbb{R}^m and set $V = \cap \mathcal{V}$. For every $\mathbf{x} \in \mathbb{R}^m$, we have:*

$$d(\mathbf{x}, V) \leq \frac{1}{1 - \xi(\mathcal{V})} \sum_{i=1}^k d(\mathbf{x}, V_i), \quad (16)$$

where ξ is given in Def. 3.

Proof: Fix $\mathbf{x} \in \mathbb{R}^m$ and $k \geq 2$. Recall that the orthogonal projection onto a subspace $V \subseteq \mathbb{R}^m$ is the mapping Π_V from \mathbb{R}^m to V that associates with each \mathbf{x} its unique nearest point in V :

$$\|\mathbf{x} - \Pi_V \mathbf{x}\|_2 = d(\mathbf{x}, V) := \min\{\|\mathbf{x} - \mathbf{v}\|_2 : \mathbf{v} \in V\}.$$

We begin the proof by observing that:

$$\begin{aligned} \|\mathbf{x} - \Pi_V \mathbf{x}\|_2 &\leq \|\mathbf{x} - \Pi_{V_k} \mathbf{x}\|_2 + \|\Pi_{V_k} \mathbf{x} - \Pi_{V_k} \Pi_{V_{k-1}} \mathbf{x}\|_2 \\ &\quad + \dots + \|\Pi_{V_k} \Pi_{V_{k-1}} \dots \Pi_{V_1} \mathbf{x} - \Pi_V \mathbf{x}\|_2 \\ &\leq \sum_{\ell=1}^k \|\mathbf{x} - \Pi_{V_\ell} \mathbf{x}\|_2 + \|\Pi_{V_k} \dots \Pi_{V_1} \mathbf{x} - \Pi_V \mathbf{x}\|_2, \end{aligned}$$

by the triangle inequality and the fact that the spectral norm $\|\Pi_{V_\ell}\|_2 \leq 1$ for all ℓ (since Π_{V_ℓ} are orthogonal projections).

The desired result (16) now follows by bounding the second term on the RHS using the following fact [17, Thm. 9.33]:

$$\|\Pi_{V_k} \Pi_{V_{k-1}} \dots \Pi_{V_1} \mathbf{x} - \Pi_V \mathbf{x}\|_2 \leq z \|\mathbf{x}\|_2, \quad (17)$$

for $z^2 = 1 - \prod_{\ell=1}^{k-1} (1 - z_\ell^2)$ and $z_\ell = \cos \theta(V_\ell, \cap_{s=\ell+1}^k V_s)$. Together with $\Pi_{V_\ell} \Pi_V = \Pi_V$ for all $\ell = 1, \dots, k$ and $\Pi_V^2 = \Pi_V$, this yields:

$$\begin{aligned} \|\Pi_{V_k} \dots \Pi_{V_1} \mathbf{x} - \Pi_V \mathbf{x}\|_2 &= \|(\Pi_{V_k} \dots \Pi_{V_1} - \Pi_V)(\mathbf{x} - \Pi_V \mathbf{x})\|_2 \\ &\leq z \|\mathbf{x} - \Pi_V \mathbf{x}\|_2. \end{aligned}$$

Since z may depend on the arbitrary labelling of the $V_i \in \mathcal{V}$, we substitute $\xi(\mathcal{V})$ for z to obtain (16). ■

Lemma 5. *Fix integers $k < m$, and let $\mathcal{T} = \{S_1, \dots, S_m\}$ be the set of intervals of length k in some cyclic order on $[m]$. Suppose there exists a map $\pi : \mathcal{T} \rightarrow \binom{[m]}{k}$ such that:*

$$\left| \bigcap_{j \in J} \pi(S_j) \right| \leq \left| \bigcap_{j \in J} S_j \right| \text{ for } J \in \binom{[m]}{k}. \quad (18)$$

Then, $|\pi(S_{j_1}) \cap \dots \cap \pi(S_{j_k})| = 1$ for all consecutive indices j_1, \dots, j_k in the order on $[m]$.

Proof: Consider the set $Q_m = \{(r, \ell) : r \in \pi(S_\ell), \ell \in [m]\}$, which has mk elements. By the pigeon-hole principle, there is some $q \in [m]$ and $J \in \binom{[m]}{k}$ such that $(q, j) \in Q_m$ for all $j \in J$. In particular, we have $q \in \cap_{j \in J} \pi(S_j)$, so that from (18) there must be some $p \in [m]$ with $p \in \cap_{j \in J} S_j$. Since $|J| = k$, this is only possible if the elements of $J = \{j_1, \dots, j_k\}$ are consecutive modulo m , in which case $|\cap_{j \in J} S_j| = 1$. Hence, it follows that $|\cap_{j \in J} \pi(S_j)| = 1$ as well.

We next consider if any other $\ell \notin J$ is such that $q \in \pi(S_\ell)$. Suppose there were such a \mathcal{T} ; then, we have $q \in \pi(S_\ell) \cap \pi(S_{j_1}) \cap \dots \cap \pi(S_{j_k})$ and (18) would imply that the intersection of every k -element subset of $\{S_\ell\} \cup \{S_j : j \in J\}$ is nonempty. This would only be possible if $\{\ell\} \cup J = [m]$, in which case the result then trivially holds since then $q \in \pi(S_j)$ for all $j \in [m]$. Suppose now there exists no such \mathcal{T} ; then, letting $Q_{m-1} \subset Q_m$ be the set of elements of Q_m not having q as a first coordinate, we have $|Q_{m-1}| = (m-1)k$.

By iterating the above arguments, we arrive at a partitioning of Q_m into sets $R_i = Q_i \setminus Q_{i-1}$ for $i = 1, \dots, m$, each having a unique element of $[m]$ as a first coordinate common to all k elements while having second coordinates which form a consecutive set. In fact, every set of k consecutive integers in the order is the set of second coordinates of some R_i . This must be the case because for every consecutive set J we have $|\cap_{j \in J} S_j| = 1$, whereas if J is the set of second coordinates for two distinct sets R_i , we would have $|\cap_{j \in J} \pi(S_j)| > 1$, violating (18). ■

Proof of Lem. 1 (Main Lemma): We assume $k \geq 2$ since the case $k = 1$ was proven in Sec. III. For simplicity, we also assume without loss of generality that the cyclic order on $[m]$ is the trivial one so that $\mathcal{T} = \{S_1, \dots, S_m\}$ has $S_i = \{i, \dots, i+k-1\}$ for $i \in [m]$ as the intervals of length k in the order. We begin by showing that $\dim(\text{Span}\{B_{\pi(S)}\}) = k$ for all $S \in \mathcal{T}$. Fix $S \in \mathcal{T}$ and note that by (8), all unit vectors $\mathbf{u} \in \text{Span}\{A_S\}$ satisfy $d(\mathbf{u}, \text{Span}\{B_{\pi(S)}\}) \leq \frac{\phi_{\mathcal{T}}(A)}{\rho^k} \delta$ for $\delta < \frac{L_2(A)}{\sqrt{2}}$. By definition of $L_2(A)$, for all 2-sparse $\mathbf{x} \in \mathbb{R}^m$:

$$L_2(A) \leq \frac{|A\mathbf{x}|_2}{\|\mathbf{x}\|_2} \leq \rho \frac{\|\mathbf{x}\|_1}{\|\mathbf{x}\|_2} \leq \rho\sqrt{2}.$$

It follows that $\delta < \rho$. Since $\phi_{\mathcal{T}}(A) \leq 1$, we also have $d(\mathbf{u}, \text{Span}\{B_{\pi(S)}\}) < 1$, and so Lem. 2 implies that $\dim(\text{Span}\{B_{\pi(S)}\}) \geq \dim(\text{Span}\{A_S\}) = k$. Since $|\pi(S)| = k$, we in fact have $\dim(\text{Span}\{B_{\pi(S)}\}) = k$, i.e. the columns of $B_{\pi(S)}$ are linearly independent.

We will now show that:

$$|\bigcap_{j \in J} \pi(S_j)| \leq |\bigcap_{j \in J} S_j|, \text{ for } J \in \binom{[m]}{k}. \quad (19)$$

Fix $J \in \binom{[m]}{k}$. By (8), for all unit vectors $\mathbf{u} \in \bigcap_{j \in J} \text{Span}\{B_{\pi(S_j)}\}$, we have that $d(\mathbf{u}, A_{S_j}) \leq \frac{\phi_{\mathcal{T}}(A)}{\rho k} \delta$ for all $j \in J$, where $\delta < \frac{L_2(A)}{\sqrt{2}}$. It follows from Lem. 4 that:

$$d\left(\mathbf{u}, \bigcap_{j \in J} \text{Span}\{A_{S_j}\}\right) \leq \frac{\delta}{\rho} \left(\frac{\phi_{\mathcal{T}}(A)}{1 - \xi(\{A_{S_j} : j \in J\})} \right) \leq \frac{\delta}{\rho},$$

where the second inequality follows from Def. 3.

Now, since $\text{Span}\{B_{\bigcap_{j \in J} \pi(S_j)}\} \subseteq \bigcap_{j \in J} \text{Span}\{B_{\pi(S_j)}\}$ and (by Lem. 3) $\bigcap_{j \in J} \text{Span}\{A_{S_j}\} = \text{Span}\{A_{\bigcap_{j \in J} S_j}\}$, we have:

$$d(\mathbf{u}, A_{\bigcap_{j \in J} S_j}) \leq \frac{\delta}{\rho}, \text{ for unit } \mathbf{u} \in \text{Span}\{B_{\bigcap_{j \in J} \pi(S_j)}\}. \quad (20)$$

In particular, Lem. 2 (since $\delta/\rho < 1$) implies that $\dim(\text{Span}\{B_{\bigcap_{j \in J} \pi(S_j)}\}) \leq \dim(\text{Span}\{A_{\bigcap_{j \in J} S_j}\})$ and (19) follows from the linear independence of the columns of A_{S_j} and $B_{\pi(S_j)}$ for all $j \in [m]$.

Suppose now that $J = \{\ell - k + 1, \dots, \ell\}$ for some $\ell \in [m]$ so that $\bigcap_{j \in J} S_j = \{\ell\}$. By (19), we have that $\bigcap_{j \in J} \pi(S_j)$ is either empty or it contains a single element. Lem. 5 ensures that the latter case is the only possibility. Thus, the association $\ell \mapsto \bigcap_{j \in J} \pi(S_j)$ defines a map $\hat{\pi} : [m] \rightarrow [m]$. Recalling (10), it follows from (20) that for all unit vectors $\mathbf{u} \in \text{Span}\{A_{\ell}\}$, we have $d(\mathbf{u}, B_{\hat{\pi}(\ell)}) \leq \delta/\rho$ also. Since ℓ is arbitrary, it follows that for every basis vector $\mathbf{e}_{\ell} \in \mathbb{R}^m$, letting $c_{\ell} = |\mathbf{A}\mathbf{e}_{\ell}|_2^{-1}$ and $\varepsilon = \delta/\rho$, there exists some $c'_{\ell} \in \mathbb{R}$ such that $|c_{\ell}\mathbf{A}\mathbf{e}_{\ell} - c'_{\ell}\mathbf{B}\mathbf{e}_{\hat{\pi}(\ell)}|_2 \leq \varepsilon$ where $\varepsilon < \frac{L_2(A)}{\sqrt{2}} \min_{\ell \in [m]} c_{\ell}$. This is exactly the supposition in (12), and the result follows from the subsequent arguments of Sec. III. ■

The strategy above can be easily modified to prove the following variation of Lem. 1, key to proving Thm. 3. Let $\pi(\mathcal{T})$ denote the set $\{\pi(S) : S \in \mathcal{T}\}$.

Lemma 6 (Main Lemma for $m < m'$). *Fix positive integers n, m, m' and k , where $k < m < m'$, and let \mathcal{T} be the set of intervals of length k in some cyclic order on $[m]$. Let $A \in \mathbb{R}^{n \times m}$ and $B \in \mathbb{R}^{n \times m'}$ both satisfy spark condition (2) with A having maximum column ℓ_2 -norm ρ . If there exists a map $\pi : \mathcal{T} \rightarrow \binom{[m']}{k}$ and some $\delta < \frac{L_2(A)}{\sqrt{2}}$ such that for $S \in \mathcal{T}$:*

$$\Theta(A_S, B_{\pi(S)}) \leq \frac{\delta}{\rho k} \min(\phi_{\mathcal{T}}(A), \phi_{\pi(\mathcal{T})}(B)), \quad (21)$$

then (9) holds for some $n \times m$ submatrix of B .

We state the required modifications briefly. Since $m' > m$, we may not invoke Lem. 5 (which requires $m = m'$) to show that $|\bigcap_{j \in J} \pi(S_j)| = 1$ when $J = \{\ell - k + 1, \dots, \ell\}$ for any $\ell \in [m]$. Instead, under the additional assumption that B satisfies the spark condition, we may simply swap the roles of A and B in the proof of (20) to show that $\dim(\text{Span}\{B_{\bigcap_{j \in J} \pi(S_j)}\}) = \dim(\text{Span}\{A_{\bigcap_{j \in J} S_j}\}) = 1$, from which the required fact then follows. The map $\hat{\pi}$ is then defined similarly, only now with codomain $[m']$, thereby reducing the proof to the $k = 1$ case

where the $n \times m$ submatrix of B is formed from the columns indexed by the image of $\hat{\pi}$.

APPENDIX B

PROOFS OF THMS. 2 & 3 AND CORS. 1 & 2

Proof of Cor. 1: We need only demonstrate how to produce N vectors \mathbf{a}_i such that for every interval of length k in some cyclic order on $[m]$, there are $(k-1)\binom{m}{k} + 1$ vectors in general linear position supported there. Let $\gamma_1, \dots, \gamma_N$ be any distinct numbers. Then the columns of the $k \times N$ matrix $V = (\gamma_i^{\ell})_{\ell, i=1}^{k, N}$ are in general linear position (since the γ_i are distinct, any $k \times k$ “Vandermonde” sub-determinant is nonzero). Next, fix a cyclic order on $[m]$ and let \mathcal{T} be the set of intervals of length k in the order. Finally, form the k -sparse vectors $\mathbf{a}_1, \dots, \mathbf{a}_N \in \mathbb{R}^m$ with supports $S \in \mathcal{T}$ (partitioning the a_i evenly among these supports so that each contains $(k-1)\binom{m}{k} + 1$ vectors \mathbf{a}_i) by setting the nonzero values \mathbf{a}_i to be those contained in the i th column of V . ■

We now determine classes of datasets Y having a stable sparse coding that are cut out by a single polynomial equation.

Proof of Thm. 2: We sketch the argument, leaving the details to the reader. Let M be the $n \times m$ matrix with columns $\mathbf{A}\mathbf{a}_i$, $i \in [N]$. Consider the following polynomial [13, Sec. IV] in the entries of A and the \mathbf{a}_i :

$$g(A, \{\mathbf{a}_i\}_{i=1}^N) = \prod_{S \in \binom{[n]}{k}} \sum_{S' \in \binom{[N]}{k}} (\det M_{S', S})^2,$$

with notation as in Sec. II.

It can be checked that when g is nonzero for a substitution of real numbers for the indeterminates, all of the genericity requirements on A and \mathbf{a}_i in our proofs of stability in Thm. 1 are satisfied (in particular, the spark condition on A). The statement of the theorem now follows directly. ■

Proof of Cor. 2: First, note that if a set of measure spaces $\{(X_{\ell}, \Sigma_{\ell}, \nu_{\ell})\}_{\ell=1}^p$ is such that ν_{ℓ} is absolutely continuous with respect to μ for all $\ell = 1, \dots, p$, where μ is the standard Borel measure on \mathbb{R} , then the product measure $\prod_{\ell=1}^p \nu_{\ell}$ is absolutely continuous with respect to the standard Borel product measure on \mathbb{R}^p (e.g., [34]). By Thm. 2, there is a polynomial that is nonzero whenever Y has a stable k -sparse representation in \mathbb{R}^m ; in particular, this property (stability) holds with probability one. ■

Proof of Thm. 3: The proof is very similar to the proof of Thm. 1 in Sec. III, the main difference being that now we establish a map $\pi : [m] \rightarrow [m']$ satisfying the requirements of Lem. 6 by pigeonholing $(k-1)\binom{m'}{k} + 1$ vectors with respect to holes $[m']$ and eventually applying Lem. 6 in place of Lem. 1. The value of C in this case is then:

$$C = \left(\frac{\sqrt{k^3}}{\min(\phi_{\mathcal{T}}(A), \phi_{\binom{[m']}{k}}(B))} \right) \frac{\max_{j \in [m]} |A_j|_2}{\min_{S \in \mathcal{T}} L_k(A X_{I(S)})}. \quad (22)$$

■